

# Stepping into the Boardroom: A novel AI-enabled framework for recognising empirical manifestations of Collective Leadership from textual data

Anonymous ACL submission

## Abstract

The concept of Collective Leadership (CL, broadly speaking leadership within groups) is difficult to define and detect empirically. A promising avenue for detecting CL focuses on discursive approaches based on group interaction and ‘turning points’ in the discussion, where participants concur on the need for action. In the absence of a defined NLP task for the detection of CL, we present a novel AI-enabled pipeline applied to publicly available hospital board text data, requiring minimal annotation thanks to in-context learning. To our knowledge, this research is the first to combine NLP and leadership theories. After presenting a language model architecture, we propose an experimental approach using ablation analysis and posit an evaluation set-up including a ‘human in the loop’ to aid acceptability by organisational research scholars and support the development of an annotated dataset.

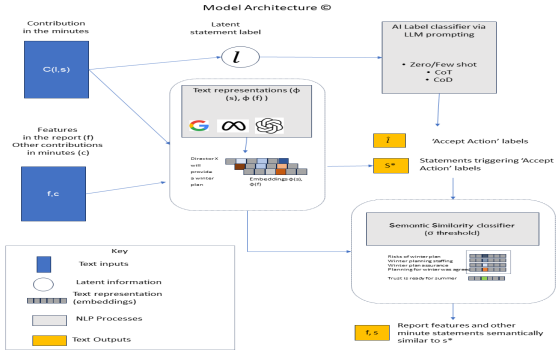


Figure 1: Proposed NLP architecture. The model receives text from the minutes and reports to create embeddings. In the first instance, a text-classification tool processes the minutes to identify formal actions (‘Accept Action’ label). Sections of the minutes classified with that label go through a semantic similarity classifier to identify other similar texts in future board reports/minutes.

## 1 Introduction and related work

The literature surrounding Collective Leadership includes ample theorising but limited research on how it manifests empirically, let alone in the context of executive boards (Edwards and Bolden, 2023; Croft et al., 2022; Ospina et al., 2020; Fairhurst et al., 2020).

Croft et al. (2022) define CL as:

“The interaction of strategic ambiguity and inward- and outward-facing reification practices to maintain divergent perspectives alongside agreed collective aims, alignment, coordination of activities, and commitment to collective success.”

A promising avenue for detecting CL and connected concepts in the above definition (such as strategic ambiguity, reification and collective work) includes discursive approaches to leadership, interaction and ‘turning points’ in a discussion, where participants concur on the need for action

(Fairhurst, 2007; Sklaveniti, 2020; Lortie et al., 2022). These discursive approaches have yet to make use of Natural Language Processing (NLP) techniques to detect CL. Our review of the NLP literature on group decision-making Mayfield and Black (2019b,a, 2020) identified only one dataset, the Wikipedia’s Article for Deletion forums (Xiao and Sitaula, 2018; Xiao and Nickerson), and no definition of an NLP task specific for detecting CL. Overall, these findings reflect that NLP (or large-scale text analytics) is hardly applied in the domain of organisational research or leadership studies (Hannigan et al., 2019). Against this background, in this study we seek to respond to this research question: *“In the absence of a defined NLP task for the detection of CL, what is the most appropriate, AI-enabled pipeline for identifying CL using solely board meeting textual data (board reports, minutes)?”*

## 2 Methods: A novel NLP task to identify CL

### 2.1 Preliminaries: a conceptual methodology for identifying CL in board text data

To inform our NLP approach, we translate the concept of CL into an executive board space of NHS hospitals by focusing on particular ‘turning points’ in the discussion, where participants formally agree on the need for change through a minuted action. From an NLP perspective, we initially seek to identify or classify sections of minutes which the language model can label as ‘Accept Action’: not only an action has been requested, but it has been formally allocated to an individual and recorded as such in the minutes. We do this based on an adaptation of the ISO standard for dialogue act annotation (ISO 24617-2:2012).

Following Croft et al.’s (2022) definition above, we posit that to detect collective leadership from executive board text requires the fulfilment of two conditions:

1. **Collective work** (joint understanding over time across teams) (Gronn, 2000): This means members of the board actively discuss an issue already highlighted in a report for that meeting, and references to an issue are seen over time (i.e. across several board sessions/reports, both in past and future). This signals a level of sustained, synergistic understanding and coordination between managers, executives and non-executives within NHS Boards.
2. **Evidence of reification over strategic ambiguity** (commitment and focus towards an aim): As noted above only those discussion points where there is reification in the form of a formally recorded action (a latent ‘Accept Action’ label) can signal CL, as these noted actions formally task managers and executives to prioritise their activities over other conflicting demands Croft et al. (2022). Consistent with the point above on collective work, we add there must be a follow-through on that specific action in future.

These degrees of reification and the distinction between collective work and collective leadership are illustrated in Figure 2. We translate these degrees of reification into a hierarchical taxonomy or ranking of ‘discussion labels’. Our analysis focuses

on the discussion label ‘Accept Action’ - a minuted action allocated to an individual. This label is more formally introduced in the next section to formulate CL mathematically, drawing from Mayfield and Black’s (2019a) notation.

### 2.2 Proposed approach

From the point of view of the meeting space, we can identify CL when we see an ‘Accept Action’ label as part of a discussion in the minutes, provided that features of that discussion will have some follow-up over time (in future), and there has been some discussion about it (contemporaneously or in the past). This is formalised in equation 2.1 below.

$$\begin{aligned} \exists \tau_1, \tau_2, \in \{1, \dots, \tau_1, \dots, t, \dots, \tau_2, \dots, T\}, \\ l_i^t = \text{‘Accept Action’} \\ \text{and } s_i^{\tau_1} \sim f_i^{\tau_1} \sim s_i^t \sim f_i^t \\ \text{and } s_i^t \sim f_i^t \sim s_i^{\tau_2} \sim f_i^{\tau_2} \end{aligned} \quad (2.1)$$

Applying a label ‘Accept Action’ at time  $t$  for a paragraph or section of the minutes  $s_i^t$  in isolation does not reflect collective leadership; it does so only if (i) we see other sections of minutes or sections of reports with a similar topic in future,  $(s_i^{\tau_2}, f_i^{\tau_2})$ , sustained commitment over time through discussion/follow-up and (ii) it has also been raised previously in minutes or reports  $(f_i^{t < \tau_1})$  - reflective of our ‘collective work’ condition above. Equation 2.2 shows there must be at least two points in time (in past  $-\tau_1$  and in future  $-\tau_2$ ) where the language model is able to identify semantic similarity compared to the text  $(s_i^t)$  in time  $t$  which has been classified with an ‘Accept Action’ label.

### 2.3 Task definition

The **input** is the dataset containing our corpus of board-level documents (reports and minutes, split in paragraphs  $f$  and  $s$  respectively) for each hospital  $h$  for the period 2017-2023. We also require a set value for the parameter sigma. In our notation  $\sigma$  is a parameter that denotes a quantitative threshold for similarity (such as Dice Coefficient or Jaccard Index (Peinelt, 2021)). As part of our experimental setting, we will test various levels of  $\sigma$ . Below we identify each report and minutes.

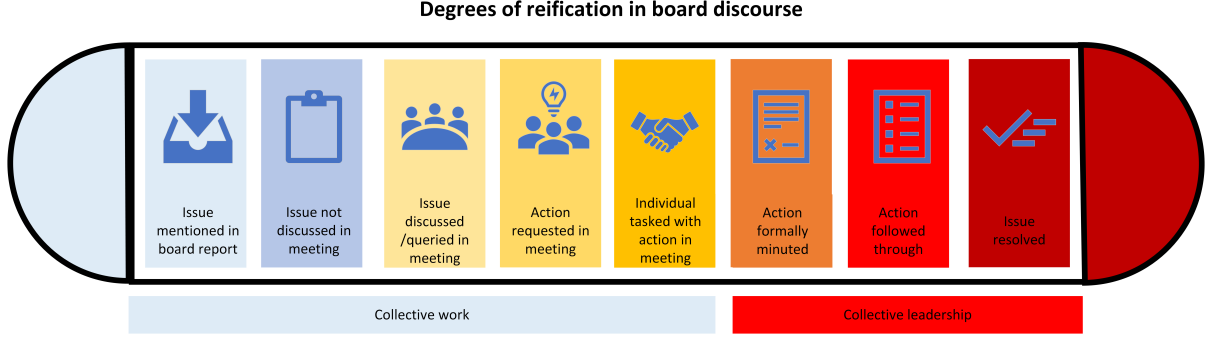


Figure 2: Degrees of reification in board discourse. The chart outlines the distinction between collective work and collective leadership

$$\begin{aligned}
 R_{ih}^t &= (l, f_1, f_2, \dots, f_F)_{ih}^t \\
 M_{ih}^t &= (c_1, c_2, \dots, c_C)_{ih}^t = \\
 &= ((l, s_1), (l_2, s_2), \dots, (l_c, s_c))_{ih}^t
 \end{aligned} \quad (2.2)$$

We defined above the text tuples  $R$  and  $C$ . These tuples are composed of a latent (unknown) action label  $l$  and a set of paragraphs (which we call features  $f$  - in reports or statements  $s$  - in the minutes).

The **objective** of this task is to identify CL as defined in equation 2.1 above for a particular  $t$  in  $\mathcal{T}$ ,  $h$  in  $\mathcal{H}$ . We do this by first identifying reification over strategic ambiguity:

$$\begin{aligned}
 \widehat{l}_{ih}^t &= \mathbb{E}(l_{ih}^t = \text{'Accept Action'} | \phi(s_{ih}^t)) \\
 \text{for any } c_{ih}^t &= (l_i, s_i)_{ih}^t \\
 s_i^* &= s_i^t \quad \text{where} \\
 \widehat{l}_i^t &= \mathbb{E}(l_i^t = \text{'Accept Action'} | \phi(s_i^t))
 \end{aligned} \quad (2.3)$$

Once we have identified the relevant section of the minutes dealing with a particular action, we verify the condition of collective work over time. We do this by identifying that a similar text which can be found in other minutes and reports at other points in time (in future).

$$\begin{aligned}
 \sigma_{ij} &= \text{Sim}(\phi(s_i^*, s_j^t)) > \sigma \quad \text{and} \\
 \varsigma_{ij} &= \text{Sim}(\phi(s_i^*, f_j^t)) > \sigma
 \end{aligned} \quad (2.4)$$

For at least one  $j$  and two  $\tau$ : one  $\tau_1 < t$  and  $\tau_2 > t$ , as per equation 2.1.

**Algorithm 1** Summary procedure for Identifying Collective Leadership from a Single Set of Reports for Hospital  $h$ , Time  $t$

**Input:** Corpus =  $(R_{ih}^t, M_{ih}^t)$  for all reports  $i$  in  $I$  and all  $t$  in  $T$ ; threshold  $\sigma$ ; label set  $L$ .

**Output:**  $O_h^\tau = \{s_i^*, s_j, f_j\}_h^\tau$

- 1: **return**  $s^*$   $\triangleright$  Identified relevant sections with 'Accept Action' label
- 2:  $s_j.append(s_j)$  for  $\tau_1 < t$   $\triangleright$  Identified relevant minutes sections where discussions took place
- 3:  $s_j.append(s_j)$  for  $\tau_2 > t$   $\triangleright$  Identified relevant sections in future meetings where the action is followed up
- 4:  $f_j.append(f_j)$  for  $\tau_1 < t$  and  $\tau_2 > t$   $\triangleright$  Identified relevant paragraphs in the minutes and reports with similar semantic similarity to each element of  $s_x$
- 5: **return**  $O_\tau$

## 2.4 Model architecture and experimental design

Figure 1 outlines the proposed architecture of the language model. Within the architecture, we have considered various potential text representations and prompting approaches as part of our experimental design.

Our experimental design considers 12 (3x4x3) architectures as outlined below:

- **Training Dataset:** We will train the classification model using in-context learning through a small manually labelled dataset drawn from minutes from a single NHS hospital not included in our sample, splitting the dataset in 80\10\10 proportion. We will aim to have at least 10 examples of each label as per Brown et al. (Brown et al., 2020), while aware the

189	source data is heavily skewed towards 'discussion' and 'query' labels instead of 'action' labels.	236
190		237
191		238
192	• <b>Text representations:</b> We will consider three different text representations using GPT-4 (by Open AI)(OpenAI, 2023), LLaMa 2 (by Meta)(noa) and BERT (by Google) (Devlin et al., 2019).	239
193		240
194		241
195		242
196		243
197	• <b>Prompts:</b> We will consider four different prompting methods: zero/few-shot (Brownlee, 2018), chain of thought prompting (Wei et al., 2023), chain of density (Adams et al., 2023).	244
198		245
199		246
200		247
201		248
202	• <b>Semantic Similarity.</b> As this is a standard NLP task, we propose to use a single architecture, tBERT (Peinelt, 2021), but testing 3 thresholds for similarity.	249
203		250
204		251
205		252
206	<b>2.5 Evaluation</b>	253
207	When developing a CL-detection NLP task we face an evaluation challenge as there is no established 'ground truth', something to benchmark the model against. In this case, the NLP literature suggests a combination of quantitative, qualitative and human-based evaluation techniques (Mayfield and Black, 2019a), which we apply to the various components of the task as well as the task overall.	254
208		255
209		256
210		257
211		258
212		259
213		260
214		261
215	Below we propose our evaluation framework, with the evaluation following in subsequent work.	262
216		263
217	• <b>Quantitative (multi-class classification):</b> Following (Brown et al., 2020), for zero/few-shot learning, we evaluate the consistency of the classification through random five-fold cross-validation against the small training dataset we have developed.	264
218		265
219		266
220		267
221		268
222		269
223	Our main evaluation metric is Balanced Accuracy, which is more sensitive to smaller class sizes. This is helpful as we have seen from preliminary review of data that the 'Accept Action' labels are less frequent than others. As a complementary metric we will consider the specific F1 for 'Accept Action' as a secondary metric, implicitly simplifying the classification problem from multi-class to binary.	270
224		271
225		272
226		273
227		274
228		275
229		276
230		277
231		278
232	• <b>Quantitative (Semantic Similarity):</b> We will use the standard F1 metric, but mindful of the lexical overlap bias found in semantic similarity tasks we consider the 'non-obvious F1' metric introduced by Peinelt (Peinelt, 2021), as a complementary metric. We will evaluate the consistency of the classification through random five-fold cross-validation throughout the small dataset constructed against Spacy's dependent parser model which has been found to have satisfactory performance in unsupervised settings.	279
233		280
234		281
235		282
		283
		284
		285
		286
		287
		288
		289
		290
		291
		292
		293
		294
		295
		296
		297
		298
		299
		300
		301
		302
		303
		304
		305
		306
		307
		308
		309
		310
		311
		312
		313
		314
		315
		316
		317
		318
		319
		320
		321
		322
		323
		324
		325
		326
		327
		328
		329
		330
		331
		332
		333
		334
		335
		336
		337
		338
		339
		340
		341
		342
		343
		344
		345
		346
		347
		348
		349
		350
		351
		352
		353
		354
		355
		356
		357
		358
		359
		360
		361
		362
		363
		364
		365
		366
		367
		368
		369
		370
		371
		372
		373
		374
		375
		376
		377
		378
		379
		380
		381
		382
		383
		384
		385
		386
		387
		388
		389
		390
		391
		392
		393
		394
		395
		396
		397
		398
		399
		400
		401
		402
		403
		404
		405
		406
		407
		408
		409
		410
		411
		412
		413
		414
		415
		416
		417
		418
		419
		420
		421
		422
		423
		424
		425
		426
		427
		428
		429
		430
		431
		432
		433
		434
		435
		436
		437
		438
		439
		440
		441
		442
		443
		444
		445
		446
		447
		448
		449
		450
		451
		452
		453
		454
		455
		456
		457
		458
		459
		460
		461
		462
		463
		464
		465
		466
		467
		468
		469
		470
		471
		472
		473
		474
		475
		476
		477
		478
		479
		480
		481
		482
		483
		484
		485
		486
		487
		488
		489
		490
		491
		492
		493
		494
		495
		496
		497
		498
		499
		500



285	We also explore potential limitations arising	
286	from various biases and challenges with repro-	
287	ducibility and accuracy arising from these methods.	
288	1. <b>Methodological biases and errors</b> might	
289	emerge through the pre-processing (encoding)	
290	of the textual data. We seek to minimise these	
291	biases by undertaking different approaches	
292	to encoding the textual data and establish-	
293	ing clear evaluation metrics for each. Re-	
294	search has also identified that Large Language	
295	Models might be biased towards outputs that	
296	mimic frequent training examples (Jones and	
297	Steinhardt). We sought to minimise this by	
298	providing a balanced set of decision-making	
299	label training examples.	
300	2. <b>Data biases.</b> Minutes, committee documents	
301	and routine reports are classified as "reporta-	
302	tive" (Heller, 2023) sources containing factual,	
303	historical information, but we recognise limi-	
304	tations related to 'authorship, bias and power'	
305	used in those documents (Heller, 2023). Large	
306	Language Models, as repositories of language	
307	data, include social biases around gender, race,	
308	religion and social constructs (Liang et al.,	
309	2021).	
310	3. <b>Researcher bias.</b> Critical to any research	
311	design is that it intimately reflects the re-	
312	searcher's perspective, which is shaped by	
313	their own beliefs and the scientific commu-	
314	nity they belong to (Kaur and Kumar, 2021).	
315	In agreement with CGT, we mitigated this by	
316	approaching the AI analysis iteratively and in	
317	a phased manner. CGT can help avoid biased	
318	interpretations of qualitative data because of	
319	this iterative approach back and forth between	
320	the human analyst and the computational anal-	
321	ysis, instead of the failed presumption that	
322	quantitative approaches are bias-free (partic-	
323	ularly given the use of natural language to	
324	'prompt' the AI) (Tschisgale et al., 2023).	
325	4. <b>Reproducibility.</b> A limitation of the approach	
326	is reproducibility, while most of the compu-	
327	tational steps are reproducible through access	
328	to the software, there is an interpreta-	
329	tion in the qualitative coding that supports	
330	grounded theory. In establishing CGT, Nel-	
331	son (Nelson, 2020) recognises that faced with	
332	the same computationally enabled results, the	
333	researcher might not code these in the same	
334	way.	
	5. <b>Accuracy of pre-trained language mod-</b>	335
	<b>els.</b> Our approach intends to build upon	336
	pre-trained large language models which are	337
	domain-agnostic. While pre-trained mod-	338
	els using domain-specific, our limited pre-	339
	annotated data might not be able to achieve	340
	higher levels of accuracy and performance,	341
	given the large cost in serving the 'long tail'	342
	of other domains (Tschisgale et al., 2023). Our	343
	training is limited to the labelling of a small	344
	section of out-of-sample board reports as to	345
	achieve a handful of examples of the different	346
	types of 'discussion labels' to classify sections	347
	of the minutes.	348
	<b>Acknowledgements</b>	349
	Thank you in advance to anonymous ARR review-	350
	ers and to [Anonymous] collaborators for their feed-	351
	back.	352
	<b>References</b>	353
	Llama 2: Open Foundation and Fine-Tuned	354
	Chat Models   Research - AI at Meta.	355
	<a href="https://ai.meta.com/research/publications/llama-2-open-foundation-and-fine-tuned-chat-models/">https://ai.meta.com/research/publications/llama-2-</a>	356
	<a href="https://ai.meta.com/research/publications/llama-2-open-foundation-and-fine-tuned-chat-models/">open-foundation-and-fine-tuned-chat-models/</a> .	357
	Griffin Adams, Alexander Fabbri, Faisal Ladhak, Eric	358
	Lehman, and Noémie Elhadad. 2023. <b>From Sparse to</b>	359
	<b>Dense: GPT-4 Summarization with Chain of Density</b>	360
	<b>Prompting.</b>	361
	Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie	362
	Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind	363
	Neelakantan, Pranav Shyam, Girish Sastry, Amanda	364
	Askell, Sandhini Agarwal, Ariel Herbert-Voss,	365
	Gretchen Krueger, Tom Henighan, Rewon Child,	366
	Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu,	367
	Clemens Winter, Christopher Hesse, Mark Chen, Eric	368
	Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess,	369
	Jack Clark, Christopher Berner, Sam McCandlish,	370
	Alec Radford, Ilya Sutskever, and Dario Amodei.	371
	2020. <b>Language Models are Few-Shot Learners.</b>	372
	Jason Brownlee. 2018. A Gentle Introduction to k-fold	373
	Cross-Validation.	374
	Charlotte Croft, Gerry McGivern, Graeme Currie, Andy	375
	Lockett, and Dimitrios Spyridonidis. 2022. <b>Unified</b>	376
	<b>Divergence and the Development of Collective Lead-</b>	377
	<b>ership.</b> <i>Journal of Management Studies</i> , 59(2):460–	378
	488.	379
	Jacob Devlin, Ming-Wei Chang, Kenton Lee, and	380
	Kristina Toutanova. 2019. <b>BERT: Pre-training of</b>	381
	<b>Deep Bidirectional Transformers for Language Un-</b>	382
	<b>derstanding.</b>	383

